

# Big Data analytics with Apache Spark and Python

## 1.Introduction

### 1.1 Context

Big data analytics cannot be narrowed down to a single tool or technology. Instead, several types of tools work together to help you collect, process, cleanse, and analyze big data. Some of the major players in big data ecosystems are: Hadoop, NoSQL databases, MapReduce, YARN, Spark.

The ability to analyze more data at a faster rate can provide big benefits to an organization, allowing it to more efficiently use data to answer important questions.

Spark is an open source cluster computing framework that uses implicit data parallelism and fault tolerance to provide an interface for programming entire clusters. Spark can handle both batch and stream processing for fast computation. [1]

### 1.2 Background

Apache Spark is an open-source, distributed processing system used for big data workloads. It utilizes in-memory caching, and optimized query execution for fast analytic queries against data of any size. It provides development APIs in Java, Scala, Python and R, and supports code reuse across multiple workloads—batch processing, interactive queries, real-time analytics, machine learning, and graph processing.

Apache Spark started in 2009 as a research project at UC Berkley's AMPLab, a collaboration involving students, researchers, and faculty, focused on data-intensive application domains. The goal of Spark was to create a new framework, optimized for fast iterative processing like machine learning, and interactive data analysis, while retaining the scalability, and fault tolerance of Hadoop MapReduce. [1]

Spark provides a faster and more general data processing platform. Spark lets you run programs up to 100x faster in memory, or 10x faster on disk, than Hadoop. Last year, Spark took over Hadoop by completing the 100 TB Daytona GraySort contest 3x faster on one tenth the number of machines and it also became the fastest open source engine for sorting a petabyte. [2]

Apache Spark architecture and key components:

Spark RDD

At the heart of Apache Spark is the concept of the Resilient Distributed Dataset (RDD), a programming abstraction that represents an immutable collection of objects that can be split across a computing cluster. [3]

### Spark SQL

Spark SQL has become more and more important to the Apache Spark project. It is the interface most commonly used by today's developers when creating applications. Spark SQL is focused on the processing of structured data, using a dataframe approach borrowed from R and Python. [3]

### Spark MLlib and MLflow

Apache Spark also bundles libraries for applying machine learning and graph analysis techniques to data at scale. MLlib includes a framework for creating machine learning pipelines, allowing for easy implementation of feature extraction, selections, and transformations on any structured dataset. [3]

MLlib (Machine Learning Library)- A distributed machine learning framework on top of Spark as a result of the distributed memory-based Spark design. Spark MLlib is 9 times as faster in comparison Hadoop disk-based version of Apache mahout (before mahout gained a Spark interface).

GraphX- A distributed graph-processing framework on top of Spark which provides an inbuilt API for expressing graph computation which will model the user-defined graphs by using the Pregel abstraction API. It additionally provides an optimized runtime for this abstraction. [5]

### Structured Streaming

Structured Streaming is a high-level API that allows developers to create infinite streaming dataframes and datasets. As of Spark 3.0, Structured Streaming is the recommended way of handling streaming data within Apache Spark, superseding the earlier Spark Streaming approach. [3]

Spark Streaming enhances the core engine of Apache Spark by providing near-real-time processing capabilities, which are essential for developing streaming analytics applications. This module can ingest live data streams from multiple sources, including Apache Kafka, Apache Flume, Amazon Kinesis, or Twitter, splitting them into discrete micro-batches.

### Key Use Cases for Spark

Generally, Spark is the best solution when time is of the essence. Apache Spark can be used for a wide variety of data processing workloads, including:

Real-time processing and insight: Spark can also be used to process data close to real-time. For example, you could use Spark Streaming to read live tweets and perform sentiment analysis on them.

Machine learning: You can use Spark MLlib to train machine learning models on large data sets and then deploy those models in your applications. It has prebuilt machine learning algorithms for tasks like regression, classification, clustering, collaborative filtering, and pattern mining.

Graph processing: You can use Spark GraphX to process graph-structured data, such as social networks or road networks. [4]

## **2. State of art**

### **2.1 analysis of different solutions/applications that use for spark sql for big data**

Apache Spark has emerged as the de facto framework for big data analytics with its advanced in-memory programming model and upper-level libraries for scalable machine learning, graph analysis, streaming and structured data processing. As the next-generation engine for big data analytics, Apache Spark can alleviate key challenges of data preprocessing, iterative algorithms, interactive analytics and operational analytics among others. With Apache Spark, data can be processed through a more general directed acyclic graph (DAG) of operators using rich sets of transformations and actions. It automatically distributes the data across the cluster and parallelizes the required operations. [6]

#### **Challenges of clustering big data**

The challenges of clustering big data are characterized into three main components:

1. Volume: as the scale of the data generated by modern technologies is rising exponentially, clustering methods become computationally expensive and do not scale up to very large datasets.
2. Velocity: this refers to the rate of speed in which data is incoming to the system. Dealing with high velocity data requires the development of more dynamic clustering methods to derive useful information in real time.
3. Variety: Current data are heterogeneous and mostly unstructured, which make the issue to manage, merge and govern data extremely challenging.

Conventional clustering algorithms cannot handle the complexity of big data due to the above reasons. For example, k-means algorithm is an NP-hard, even when the number of clusters is small. Consequently, scalability is a major challenge in big data. [7]

Spark SQL is a component of Apache Spark that provides a unified interface for querying structured and semi-structured data using SQL, DataFrame, and Dataset APIs. It allows users to seamlessly integrate SQL queries with Spark's distributed processing capabilities, enabling efficient analysis of large-scale data. Here's an analysis of different solutions and applications for Spark SQL in the context of big data:

#### **1. Data Warehousing:**

- Spark SQL can be used as part of a data warehousing solution to analyze large volumes of structured data stored in data lakes or cloud storage platforms like Hadoop Distributed File System (HDFS), Amazon S3, or Azure Data Lake Storage.
- Organizations can leverage Spark SQL to perform complex analytics, run ad-hoc queries, and generate reports on data warehousing platforms built using Apache Hive, Apache Hadoop, or cloud-based data warehousing solutions.

#### **2. Interactive Data Analysis and Business Intelligence:**

- Spark SQL enables interactive data analysis and exploration by providing a SQL interface to query large datasets stored in distributed file systems or databases.
- Business intelligence tools like Tableau, Qlik, and Looker can integrate with Spark SQL to enable self-service analytics, data visualization, and reporting on big data platforms.

#### **3. Data Integration and ETL:**

- Spark SQL facilitates data integration and ETL (Extract, Transform, Load) processes by providing SQL-based access to diverse data sources such as JSON, CSV, Parquet, Avro, JDBC, and ORC.
- Solutions like Apache NiFi, Talend, and Apache Airflow integrate with Spark SQL to ingest, transform, and load data from various sources into data lakes or analytical databases for further processing and analysis.

#### **4. Structured Streaming and Real-time Analytics:**

- Spark SQL's Structured Streaming API allows users to perform real-time stream processing and analytics on continuous data streams.

- Organizations can build real-time analytics applications using Spark SQL to analyze streaming data from sources like Apache Kafka, Amazon Kinesis, or Azure Event Hubs and derive insights in near real-time.

#### **5. Machine Learning and Predictive Analytics:**

- Spark SQL can be integrated with Spark's machine learning library (MLlib) to perform SQL-based feature engineering, data preprocessing, and model training on large datasets.
- Data scientists and analysts can leverage Spark SQL to prepare data for machine learning pipelines, run SQL queries to explore data patterns, and generate features for predictive modeling tasks.

#### **6. Data Governance and Security:**

- Spark SQL provides features for data governance and security, including finegrained access control, encryption, and auditing capabilities.
- Organizations can enforce data governance policies, control access to sensitive data, and monitor SQL queries executed on big data platforms using Spark SQL.

#### **7. Data Lake Analytics:**

- Spark SQL plays a crucial role in data lake analytics, allowing users to query and analyze data stored in data lakes using SQL.
- By leveraging Spark SQL, organizations can unlock the value of data lakes by performing ad-hoc analysis, generating insights, and building data-driven applications on top of large-scale data lake repositories.

#### **8. Cross-Platform Data Analysis and Integration:**

- Spark SQL supports cross-platform data analysis and integration by providing connectors and APIs for interacting with various data sources and platforms, including traditional databases, cloud storage, NoSQL databases, and data warehouses.
- Organizations can use Spark SQL to integrate and analyze data from disparate sources, enabling unified data analytics across hybrid cloud and multi-cloud environments.

## **2.2 Comparisons**

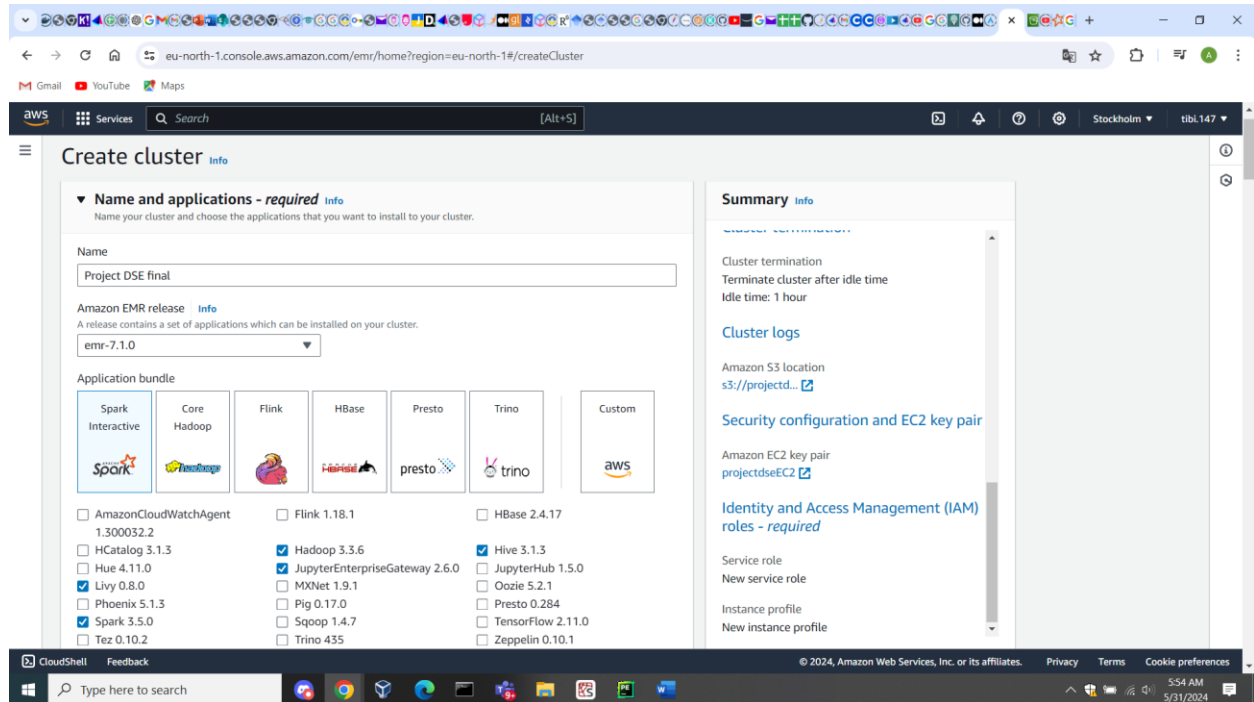
Apache Hadoop and Apache Spark are both highly capable platforms, and the features they offer cater to different data processing needs. In addition to the previously mentioned features, some other important features that differentiate these platforms include APIs for data processing, cluster management systems, machine learning libraries, graph processing libraries, real-time data processing, data visualization tools, and stream processing frameworks.

Apache Hadoop offers APIs such as MapReduce, Pig, Hive, HBase, Sqoop, Flume, and Oozie. It has cluster management systems such as Apache Ambari and Cloudera Manager, and machine learning libraries such as Mahout and MLlib. Hadoop has graph processing libraries like Giraph and Hama, and real-time data processing frameworks such as Kafka, Storm, and Samza. It also has data visualization tools like Apache Zeppelin and Hue.

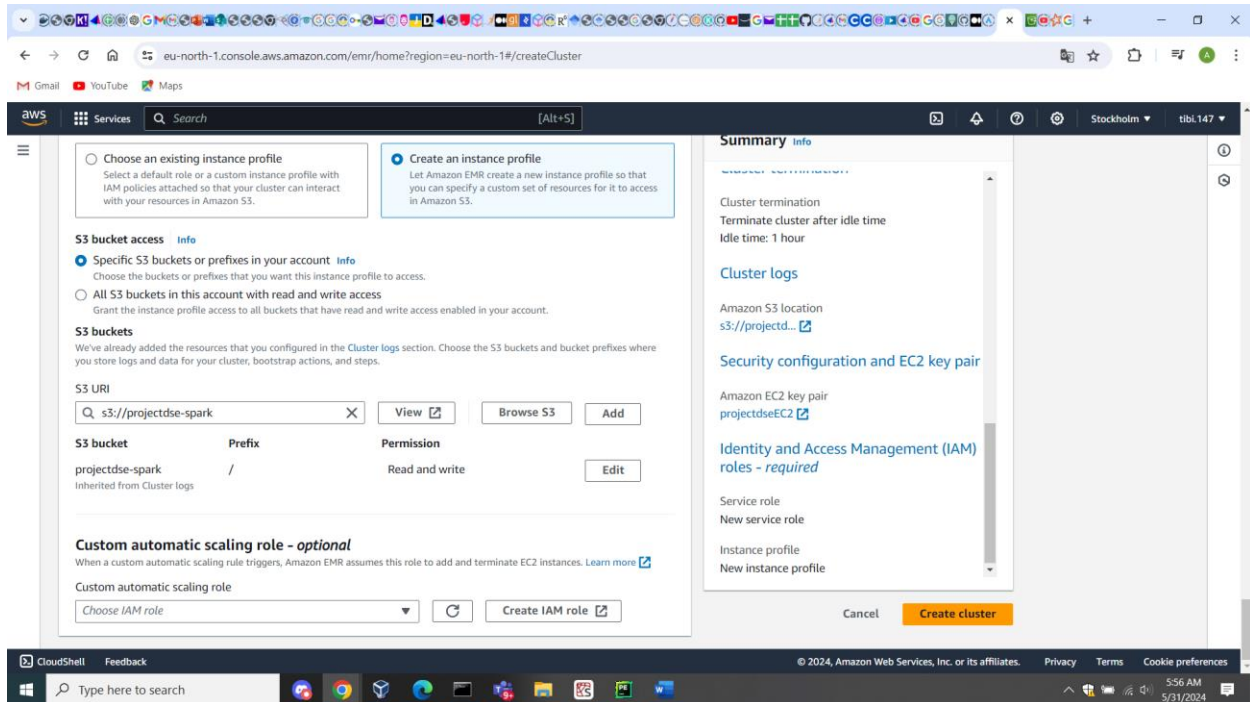
Apache Spark offers APIs like Spark SQL, Spark Streaming, Spark MLlib, and others. It has cluster management systems like Apache Mesos, Hadoop YARN, and Kubernetes, and machine learning libraries like MLlib, GraphFrames, and TensorFlow on Spark. Spark has graph processing libraries like GraphX, GraphFrames, and Neptune and real-time data processing frameworks such as Spark Streaming, Flink, and Storm. It also has data visualization tools like Apache Zeppelin, Jupyter Notebook, and others.

## **Phase 5 Final results and conclusions**

- 1.



The final solution involves setting up and configuring an Apache Spark cluster on Amazon Web Services (AWS) to efficiently process and analyze large datasets. This setup leverages AWS's Elastic MapReduce (EMR) service, which simplifies the deployment and management of Spark clusters and run Apache Spark and other big data frameworks. Amazon S3 (Simple Storage Service) it is used for storing input data, intermediate results, and final output. Amazon EC2 (Elastic Compute Cloud) the hey pair of virtual servers to host the Spark cluster nodes (Master and Workers). The solution ensures high availability, scalability, and cost-effectiveness, while also integrating with other AWS services for enhanced functionality.



## Cluster Setup:

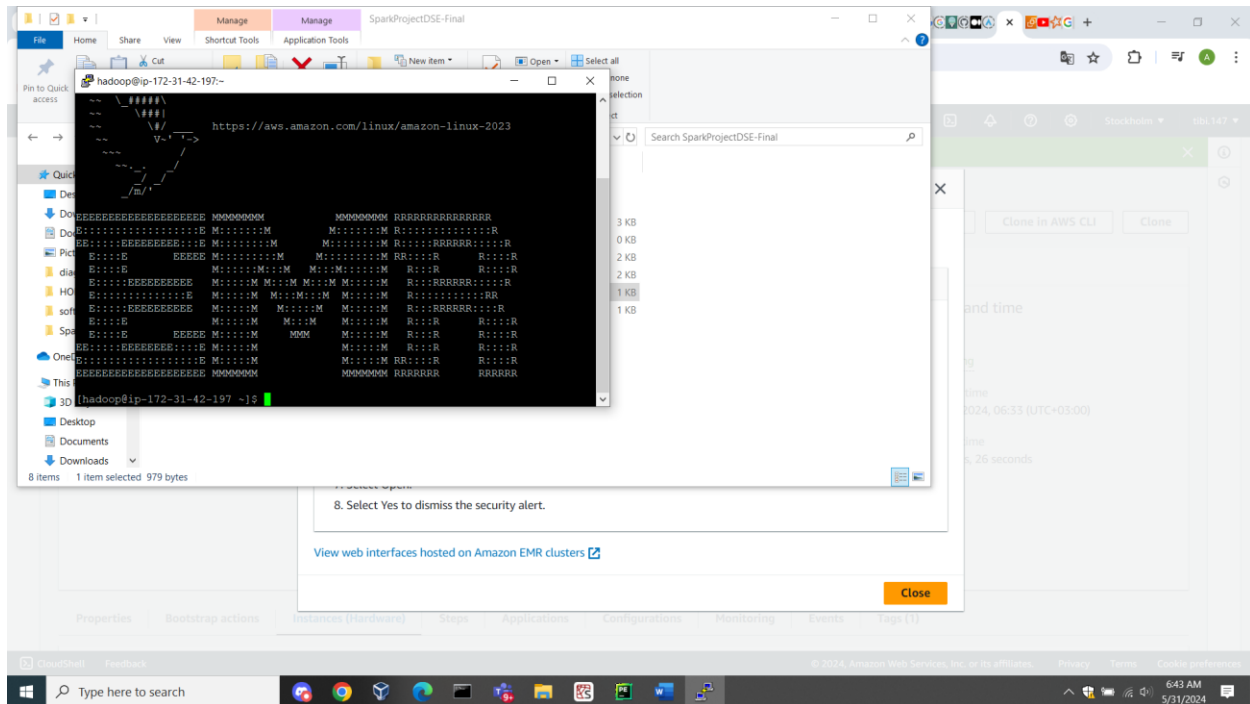
Master Node manages the cluster, coordinates tasks, and maintains cluster state and Worker Nodes execute tasks and process data. The number of worker nodes can be scaled up or down based on workload. Amazon S3 acts as the primary data lake for input, output, and intermediate data. Amazon EMR handles scaling, monitoring, and maintenance of the cluster. Auto-scaling can be configured to adjust the number of worker nodes based on the workload.

Spark reads data from Amazon S3, processes it, and writes results back to Amazon S3. Logs are stored in Amazon S3 for debugging and analysis.



2.

After I created the cluster, I opened the Elastic MapReduce (EMR) service in puTTY terminal by introducing the cluster's host name and the private key that I created in EC2 (Elastic Compute Cloud).



In the terminal, I generated two codes to copy my python script and dataset script from Amazon S3 bucket to the local file system.



3.

The results that I noticed consisted of:

Performance:

The Spark cluster on AWS can handle varying workloads efficiently by scaling the number of worker nodes up or down as needed. This elasticity ensures that performance remains optimal even as data volumes and computational demands increase. With appropriate instance types and optimized cluster configurations, Spark jobs run significantly faster compared to traditional on-premises setups. The use of high-performance EC2 instances and optimized networking enhances processing speed.

Management and Maintenance:

Using Amazon S3 as a data lake allows for efficient data ingestion, processing, and storage, enabling a unified data platform. Built-in redundancy and fault tolerance in AWS services ensure high availability and reliability of the Spark cluster.

The screenshot displays the Amazon EMR History Server web interface. At the top, the browser address bar shows the URL: `p-m6aj6h7qbt7-shs.emrappui-prod.eu-north-1.amazonaws.com/shs/?showIncomplete=true`. The page header includes the Amazon EMR logo and the text "History Server". Below the header, the "Event log directory" is listed as `s3a://prod.eu-north-1.appinfo.src/j-RN6AJ6H7QBT7/sparklogs`, with a "Last updated" timestamp of 2024-05-31 06:55:08 and a "Client local time zone" of Europe/Bucharest. A search bar is located on the right side of the page. The main content area features a table with the following columns: Version, App ID, App Name, Started, Spark User, Last Updated, and Event Log. A single entry is shown in the table, representing an application named "MovieSimilarities1M.py" that started on 2024-05-31 at 06:49:44 and was last updated on 2024-05-31 at 06:54:48. The application was run by the "hadoop" user. A "Download" button is available for the event log. Below the table, it indicates "Showing 1 to 1 of 1 entries" and provides a link to "Back to completed applications". The bottom of the screenshot shows the Windows taskbar with the search bar and various application icons.

Version	App ID	App Name	Started	Spark User	Last Updated	Event Log
3.5.0-amzn-1	application_1717126679793_0001	MovieSimilarities1M.py	2024-05-31 06:49:44	hadoop	2024-05-31 06:54:48	<a href="#">Download</a>

Observations:

Setting up a Spark cluster using Amazon EMR is straightforward, with options to customize configurations based on specific needs.

Tuning Spark configurations (executor memory, number of executors) and data partitioning strategies resulted in substantial performance improvements.

The choice of instance types (compute-optimized vs. memory-optimized) has a significant impact on performance. For memory-intensive tasks, memory-optimized instances showed better performance.

Tuning Spark configurations (executor memory, number of executors) and data partitioning strategies resulted in substantial performance improvements.

The screenshot shows the Amazon EMR console interface. At the top, there's a navigation bar with the Hadoop logo and 'All Applications' title. Below this is a table of applications. The table has columns: ID, User, Name, Application Type, Application Tags, Queue, Application Priority, StartTime, LaunchTime, FinishTime, and State. A single application is listed with ID 'application\_1717126679793\_0001', User 'hadoop', Name 'MovieSimilarities1M.py', Application Type 'SPARK', Queue 'default', Application Priority '0', StartTime 'Fri May 31 06:50:14 +0300 2024', LaunchTime 'Fri May 31 06:50:15 +0300 2024', FinishTime 'N/A', and State 'RUNNING'. The left sidebar contains 'Application History' and 'Tools' sections. The bottom of the image shows a Windows taskbar with various application icons and the system clock indicating 6:54 AM on 5/31/2024.

ID	User	Name	Application Type	Application Tags	Queue	Application Priority	StartTime	LaunchTime	FinishTime	State
application_1717126679793_0001	hadoop	MovieSimilarities1M.py	SPARK		default	0	Fri May 31 06:50:14 +0300 2024	Fri May 31 06:50:15 +0300 2024	N/A	RUNNING

Discussions:

Advantages of AWS for Spark:

The managed nature of Amazon EMR reduces the complexity of running and maintaining a Spark cluster, allowing teams to focus on data processing and analysis.

AWS's scalability ensures that the Spark cluster can grow with the data and processing needs, while the flexibility in configurations allows for tailored solutions.

## Challenges and Solutions:

While spot instances offer cost savings, their potential interruptions require robust checkpointing and job rescheduling mechanisms to ensure reliability.

Transferring large datasets in and out of AWS can incur significant costs. Utilizing AWS Direct Connect or optimizing data transfer strategies can help mitigate these expenses.

## Future Considerations:

Integrating Spark with AWS's machine learning services can enhance data processing workflows with advanced analytics and predictive modeling capabilities.

## 4. Conclusion:

Running Apache Spark on AWS using Amazon EMR offers a powerful, flexible, and cost-effective solution for big data processing. The combination of Spark's computational capabilities and AWS's scalable infrastructure provides a robust platform for handling diverse data workloads. Implementing best practices, such as efficient data partitioning, proper instance type selection, and optimized Spark configurations, leads to significant performance gains and cost savings. These practices ensure that the Spark cluster operates at peak efficiency.

## Bibliography

- [1] "<https://www.tableau.com/learn/articles/big-data-analytics>," [Online].
- [2] "<https://aws.amazon.com/what-is/apache-spark/>," [Online].
- [3] "<https://www.toptal.com/spark/introduction-to-apache-spark>," [Online].
- [4] "<https://www.infoworld.com/article/3236869/what-is-apache-spark-the-big-data-platform-thatcrushed-hadoop.html>," [Online].
- [5] C. S. ., G. S. S. D. S. I. K. Monika Chand, "Analysis of Big Data using Apache Spark".
- [6] "<https://nexocode.com/blog/posts/what-is-apache-spark/>," [Online].
- [7] S. S. . R. D. . X. C. . P. X. P. ., "Big data analytics on Apache Spark," *Springer International Publishing Switzerland*, 2016.
- [8] A. M. R. R. K. R. K. P. Kumar, "A Review on Apache Spark".
- [9] "What is Apache Spark," [Online]. Available: <https://aws.amazon.com/what-is/apache-spark/>.
- [10] "<https://aws.amazon.com/what-is/apache-spark/>," [Online].
- [11] M. M. Saeed1, "Big data clustering techniques based on," *PeerJ Computer Science*, 2020.
- [12] "Analysis of Big Data using Apache Spark".

